# Probability and Statistics for Ensemble Forecasting

Tom Hamill (NOAA/ESRL, Boulder)
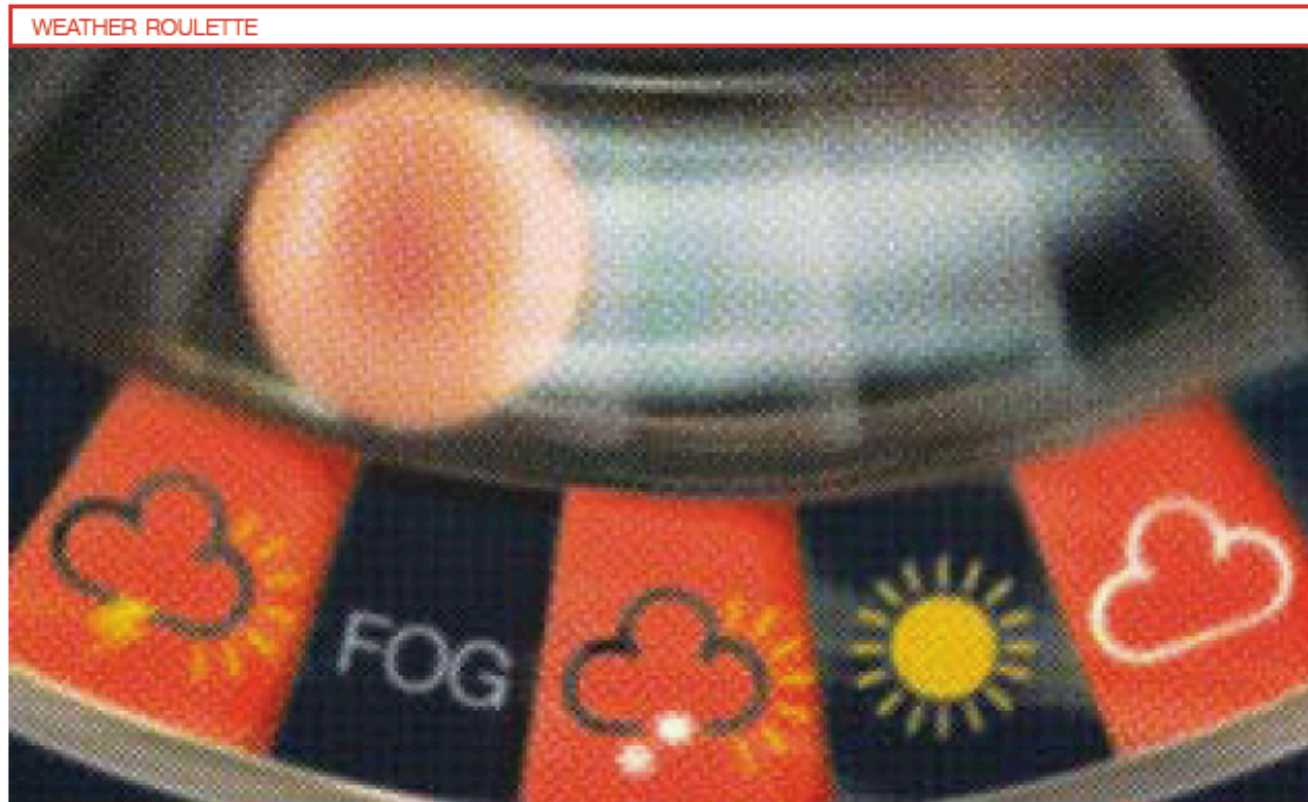and
Jim Hansen (Navy/NRL, Monterey)

(borrows heavily from Dan Wilks'
*Statistical Methods in the Atmospheric Sciences text*)

# Probability and statistics

- **Probability**: a formalism for expressing uncertainty quantitatively.

- **Statistics**: the science pertaining to the collection, analysis, interpretation or explanation, and presentation of data.

- **Goal**: get you comfortable with the terminology the other instructors will use.
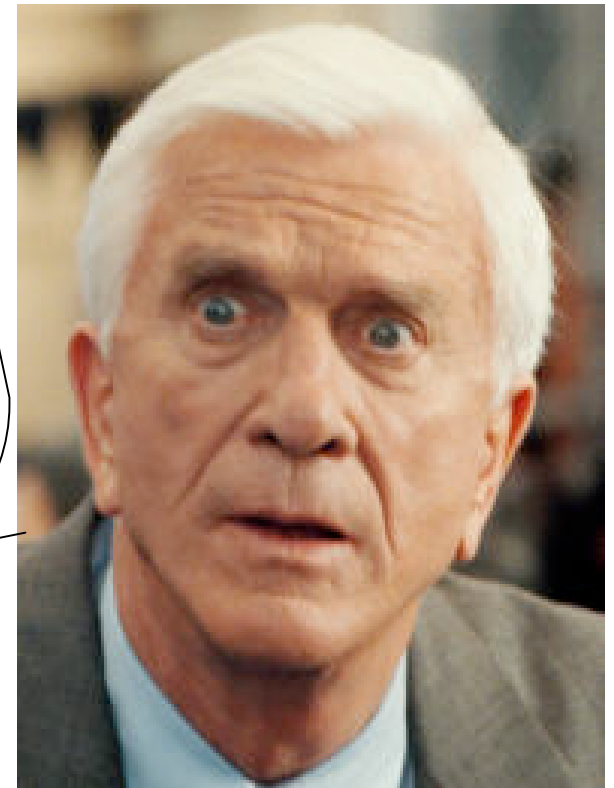
# Part 1: Probability



WEATHER ROULETTE

FOG

Weather is uncertain, so we use the language of uncertainty

photo courtesy of Lenny Smith, Oxford U. and London School of Economics

# Is probability
## (1) inherently confusing, or
## (2) a formal way of bamboozling and waffling?



"Doctors say that Nordberg has a 50/50 chance of living, though there's only a 10 percent chance of that."

# Axioms of Probability

S

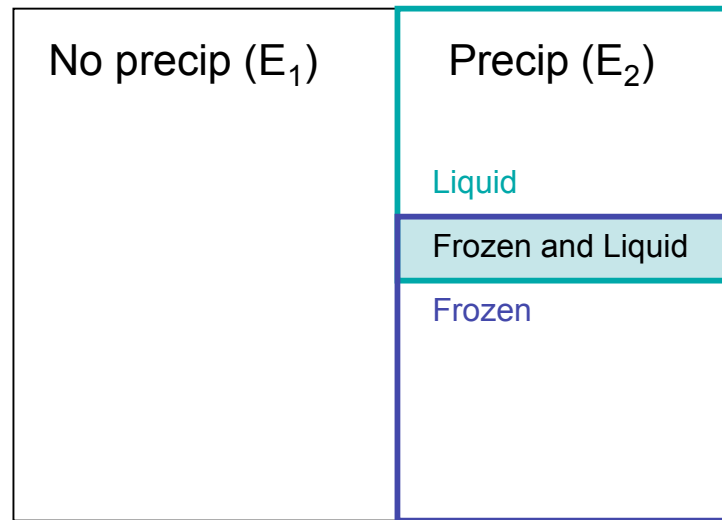| No precip ($E_1$) | Precip ($E_2$) |
|---|---|
| | Liquid |
| | Frozen and Liquid |
| | Frozen |

$0.0 \leq Pr(E_1) \leq 1.0$

$Pr(S) = 1.0$

$Pr(E_1) + Pr(E_2) = 1.0$

S is the "sample space." $E_1$ and $E_2$ are "mutually exclusive" and "collectively exhaustive" events that fill the sample space.
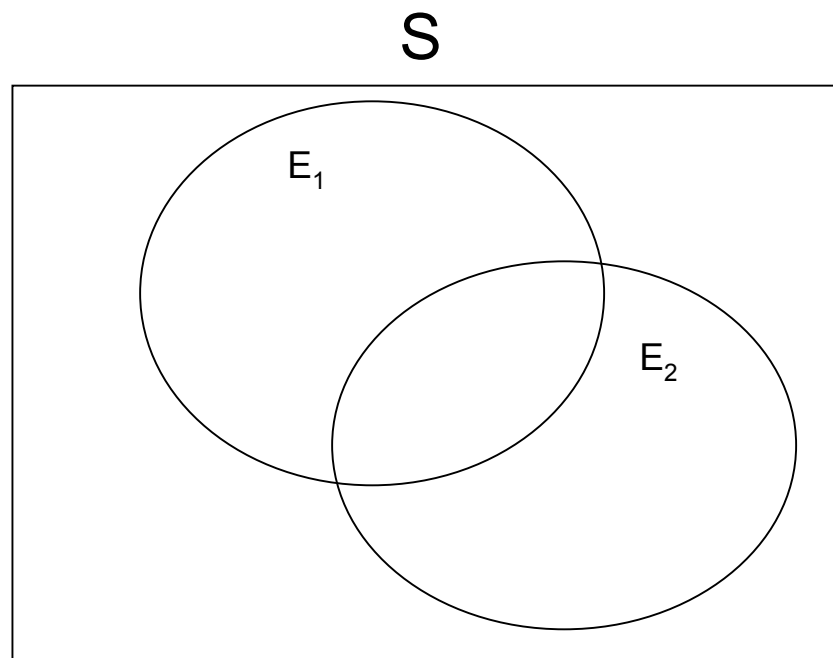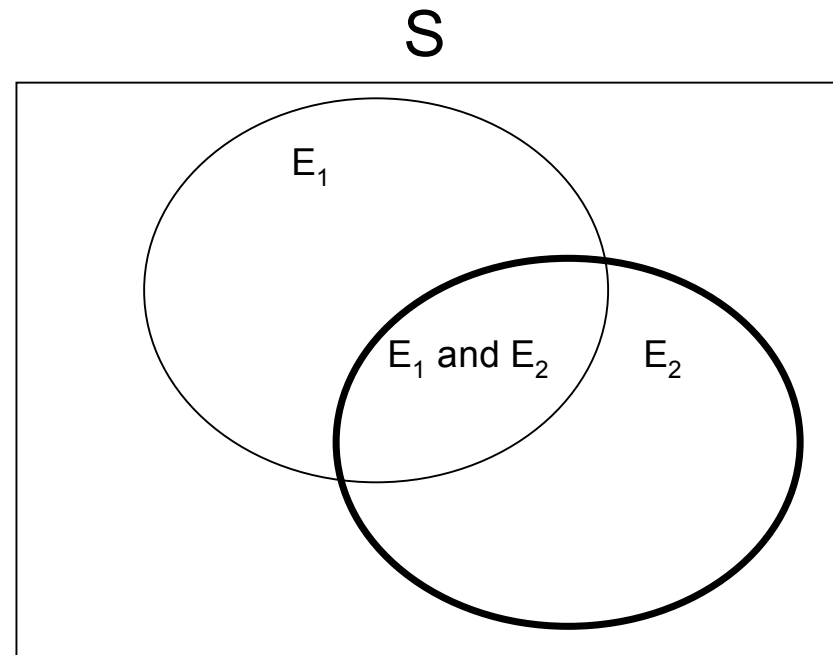
# Union of Events

S



Pr(Liquid) and Pr(Frozen) = Pr(Liquid) + Pr(Frozen) - Pr(Frozen and Liquid)

# Conditional Probability

# Conditional Probability



$\Pr(E_1 \mid E_2) = \Pr(E_1$ given that $E_2$ has occurred$)$
$= \Pr(E_1$ and $E_2) / \Pr(E_2)$

narrow the playing field … consider only the subset where $E_2$ has occurred
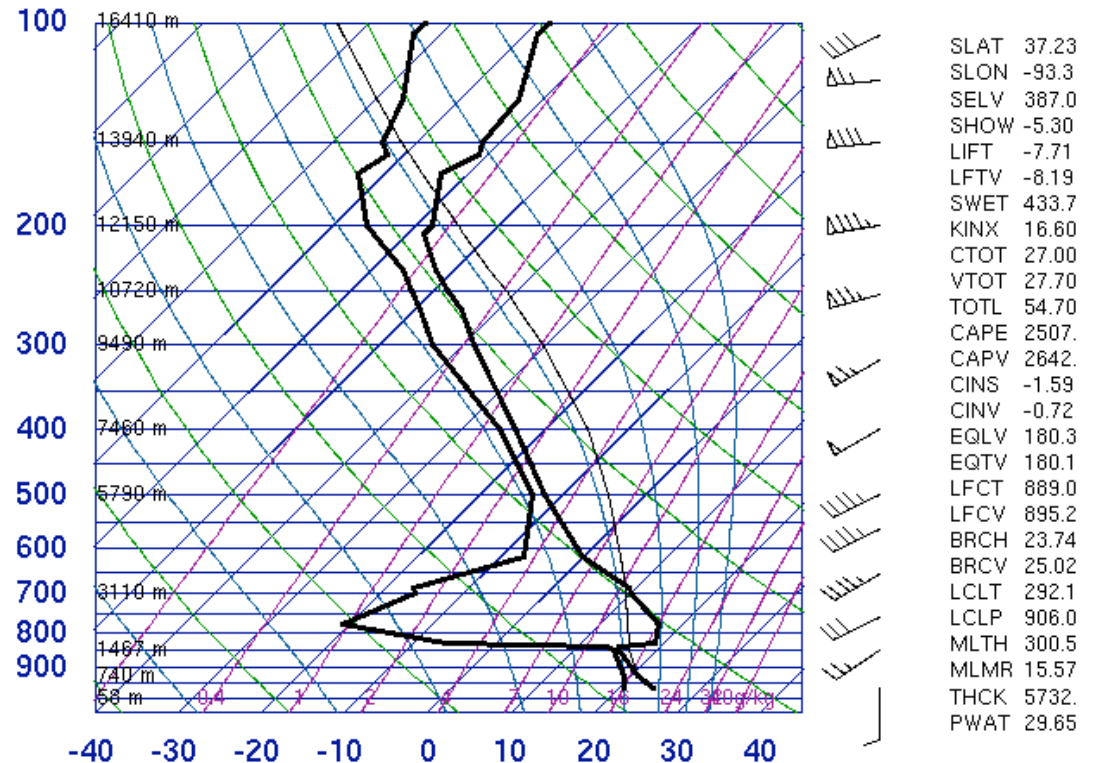
# Example: "loaded gun" sounding

P(tornado in SW MO) = 0.02

unconditional probability
of a tornado is small; most
likely it will be impossible
to break through the
capping inversion.

P(tornado | thunderstorm) = 0.35

if penetrative convection
does happen, the large
instability and shear increase
the probability that the
thunderstorm will produce a tornado.



72440 SGF Springfield

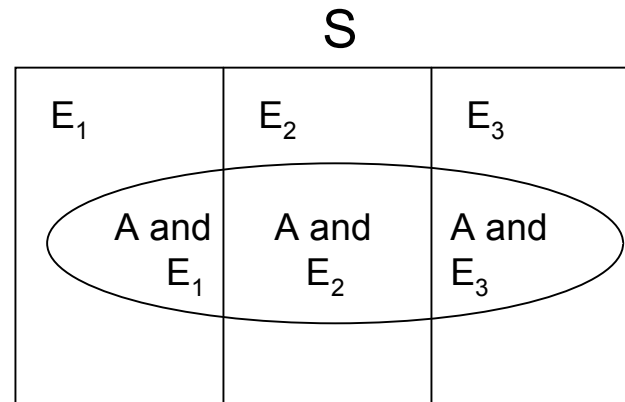| | |
|---|---|
| SLAT | 37.23 |
| SLON | -93.3 |
| SELV | 387.0 |
| SHOW | -5.30 |
| LIFT | -7.71 |
| LFTV | -8.19 |
| SWET | 433.7 |
| KINX | 16.60 |
| CTOT | 27.00 |
| VTOT | 27.70 |
| TOTL | 54.70 |
| CAPE | 2507. |
| CAPV | 2642. |
| CINS | -1.59 |
| CINV | -0.72 |
| EQLV | 180.3 |
| EQTV | 180.1 |
| LFCT | 889.0 |
| LFCV | 895.2 |
| BRCH | 23.74 |
| BRCV | 25.02 |
| LCLT | 292.1 |
| LCLP | 906.0 |
| MLTH | 300.5 |
| MLMR | 15.57 |
| THCK | 5732. |
| PWAT | 29.65 |

00Z 10 May 2003

University of Wyoming

# Independence

- $E_1$ and $E_2$ are independent if and only if $\Pr(E_1 \text{ and } E_2) = \Pr(E_1) \times \Pr(E_2)$
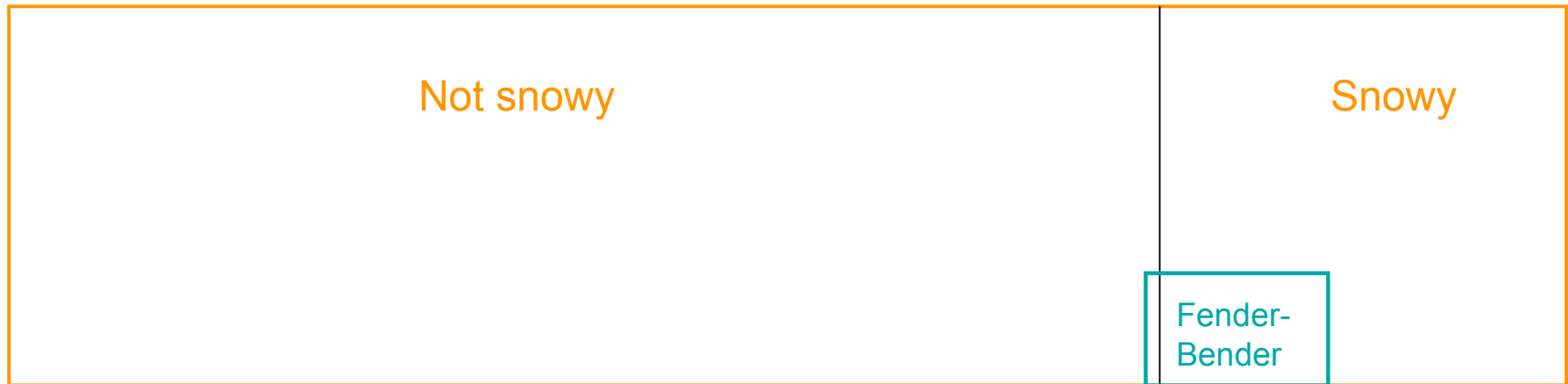


Probability of two sixes =
1/6 x 1/6 = 1/36

# Law of total probability

S

| $E_1$ | $E_2$ | $E_3$ |
|---|---|---|
| A and $E_1$ | A and $E_2$ | A and $E_3$ |

$$\Pr(A) = \sum_{i=1}^{3} \Pr(A \mid E_i) \Pr(E_i)$$

Overall "unconditional" probability can be computed summing / integrating the weighted conditional probabilities

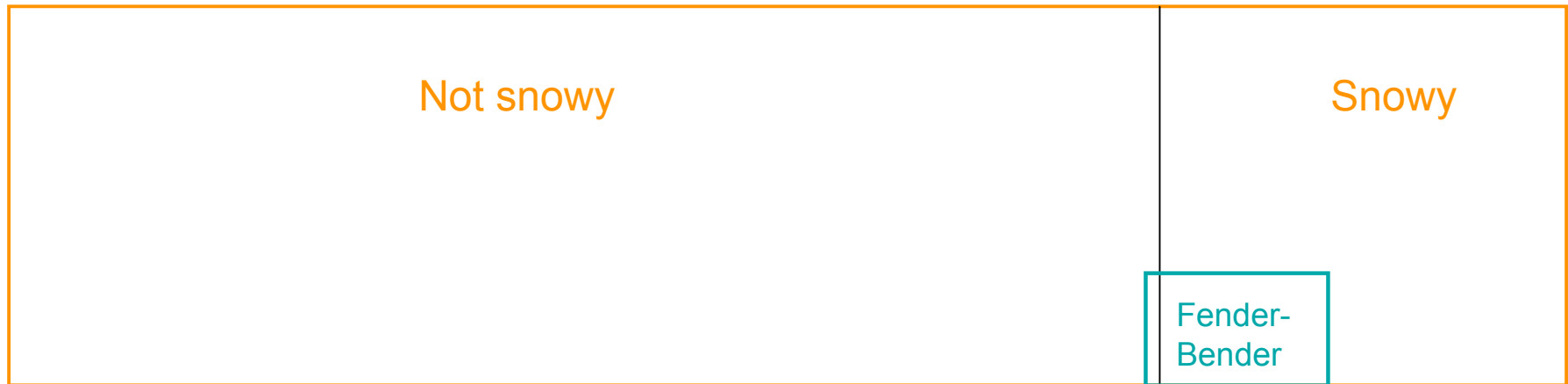# Law of total probability: driving example



P(Fender-Bender) = P(Fender-Bender| Not Snowy) P(Not Snowy)  +
P(Fender-Bender| Snowy) P(Snowy)

= 0.01         x   0.75       +
0.10         x   0.25

= 0.0325

# Law of total probability: driving example

| Not snowy | Snowy |
|---|---|
|  |  |
|  | Fender-Bender |

P(Fender-Bender) = P(Fender-Bender| Not Snowy) P(Not Snowy)  +
P(Fender-Bender| Snowy) P(Snowy)

= 0.01                    x    0.75             +
0.10                      x    0.25

= 0.0325   (I'm an excellent driver)

# Discrete vs. Continuous Probability
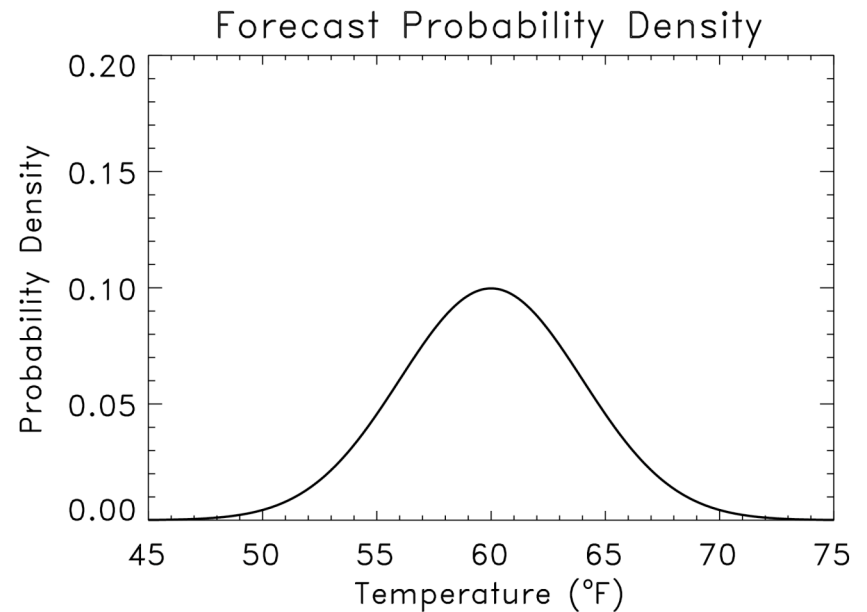
- **Discrete**: limited number of possible outcomes

- **Continuous**: unlimited number of outcomes.

P(T=60.0) not meaningful;

probability density
expressed relative likelihood
of being *near* a particular value;
and probability density follows
other probability axioms, e.g.,

$$\int_{t=0K}^{\infty} P(t)dt = 1.0$$



Forecast Probability Density

# Discrete "parametric" probability distributions: the binomial distribution

$$\Pr(X = x) = \binom{N}{x} p^x (1-p)^{N-x}$$

$$\binom{N}{x} = \frac{N!}{x!(N-x)!}$$
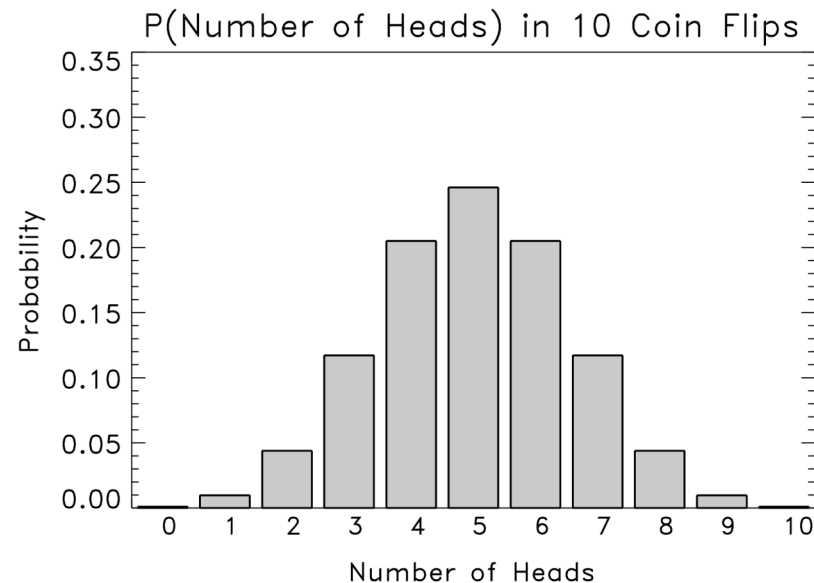
$X$ is random variable

$x$ is a specific number

$N$ is the number of trials

$p$ is the event probability

P(Number of Heads) in 10 Coin Flips

# Discrete "parametric" probability distributions: the binomial distribution

$$\Pr(X = x) = \binom{N}{x} p^x (1-p)^{N-x}$$

$$\binom{N}{x} = \frac{N!}{x!(N-x)!}$$
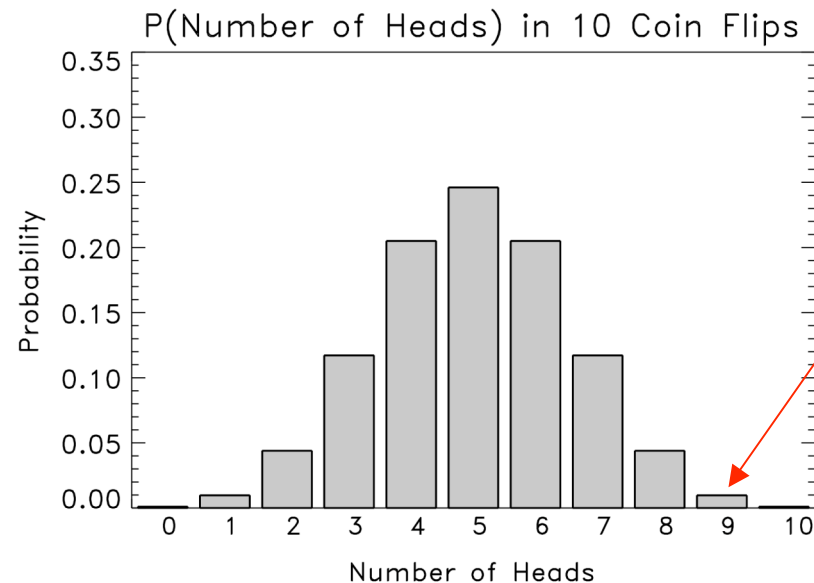
*X* is random variable

*x* is a specific number

*N* is the number of trials

*p* is the event probability

P(Number of Heads) in 10 Coin Flips
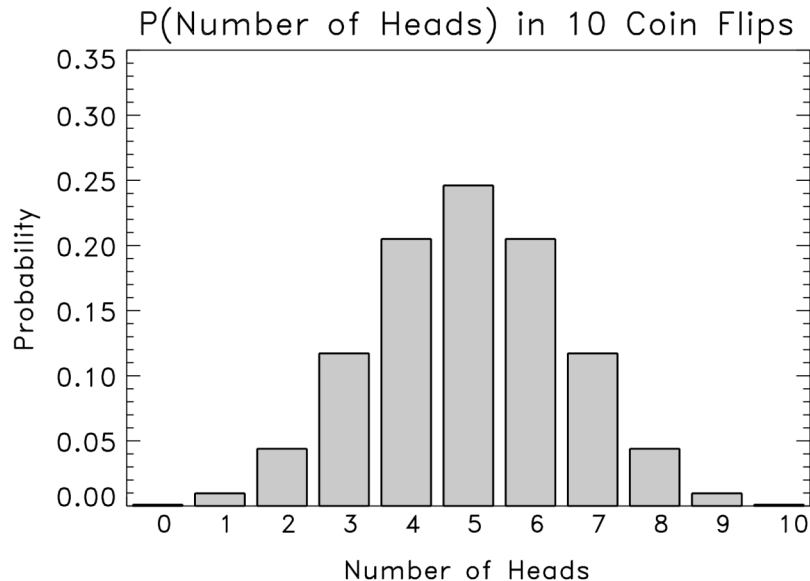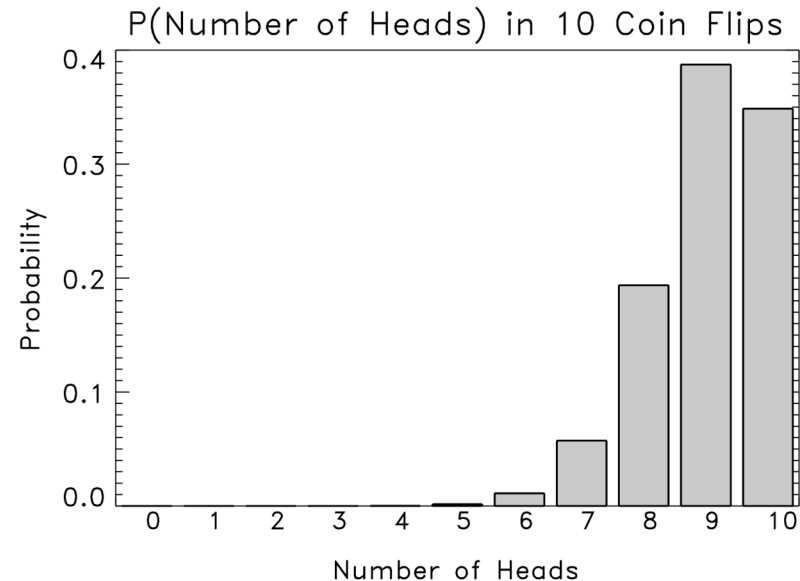
analogy: you might forecast 50% probability of rain, and rain may happen 9/10 times. That can happen, though it's unlikely.

# Binomial distributions

p = 0.5                              p = 0.9



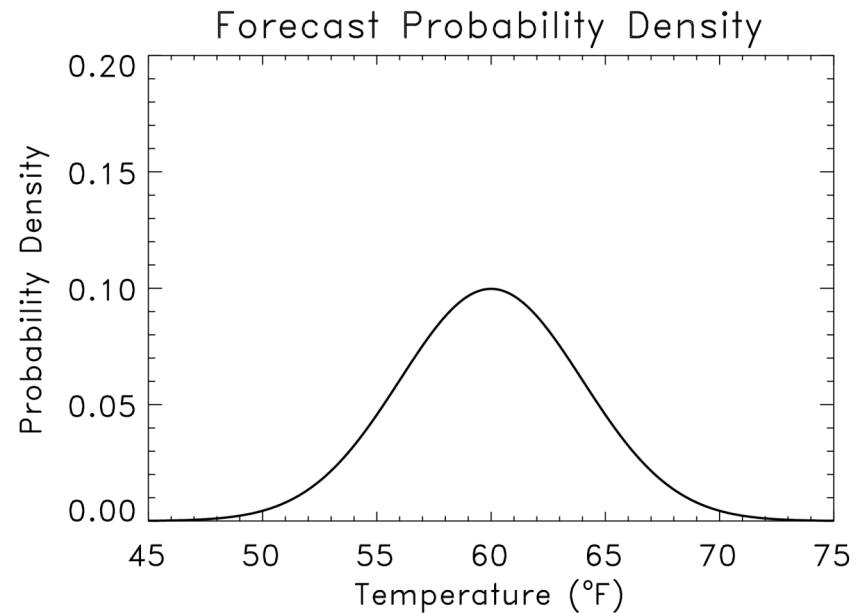…though perhaps p = 0.9 would have been a more appropriate choice.

# Continuous parametric probability distributions: the Normal distribution

- Also called "Gaussian" or 'the bell-shaped curve"

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

- $f(x)$ is the "probability density"
- $\mu$ is the "mean"
- $\sigma$ is the "standard deviation"

$\mu = 60.0, \ \sigma = 4.0$



Forecast Probability Density

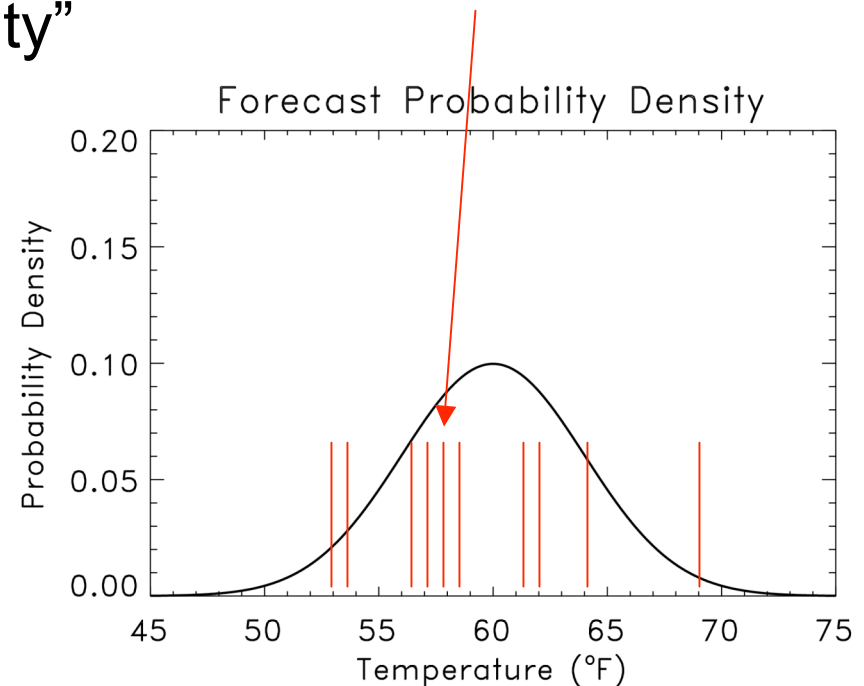# Continuous parametric probability distributions: the Normal distribution

- Also called "Gaussian" or 'the bell-shaped curve"

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

this ensemble might be a random sample from a smooth distribution like this

- $f(x)$ is the "probability density" function, or PDF
- $\mu$ is the "mean"
- $\sigma$ is the "standard deviation"
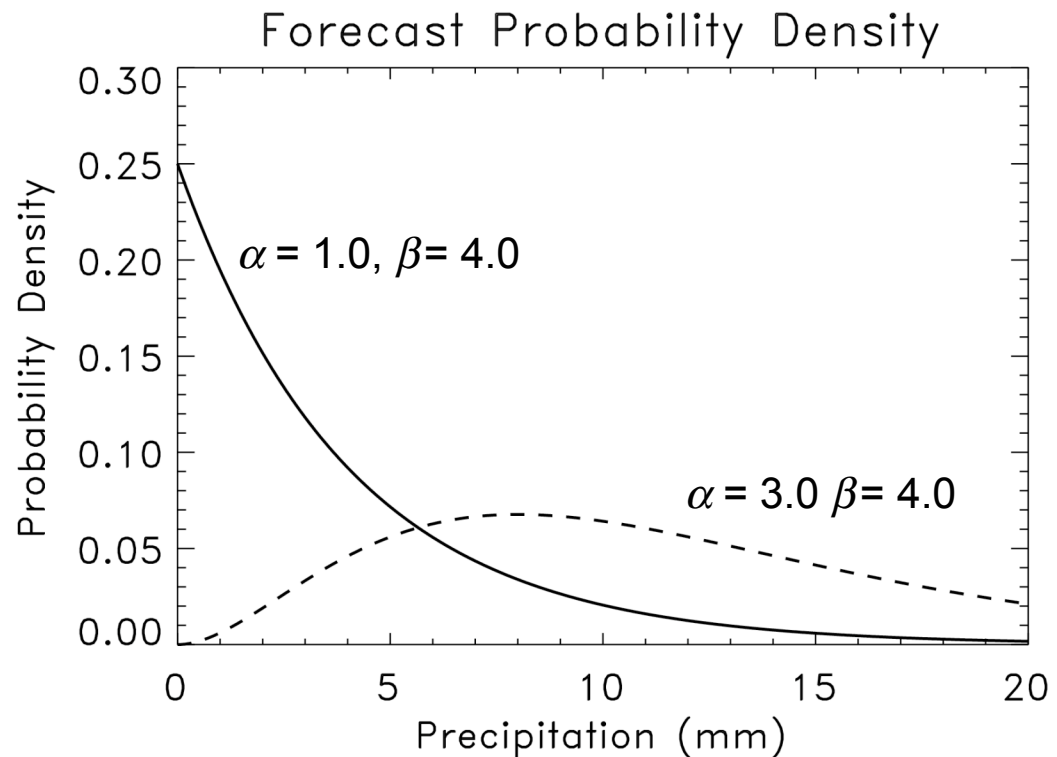
$\mu = 60.0, \ \sigma = 4.0$

Forecast Probability Density

# The gamma distribution

$$f(x) = \frac{\left(x/\beta\right)^{\alpha-1}\exp\left(-x/\beta\right)}{\beta\,\Gamma(\alpha)}, \qquad x,\alpha,\beta > 0.0$$

$\alpha$ = shape parameter
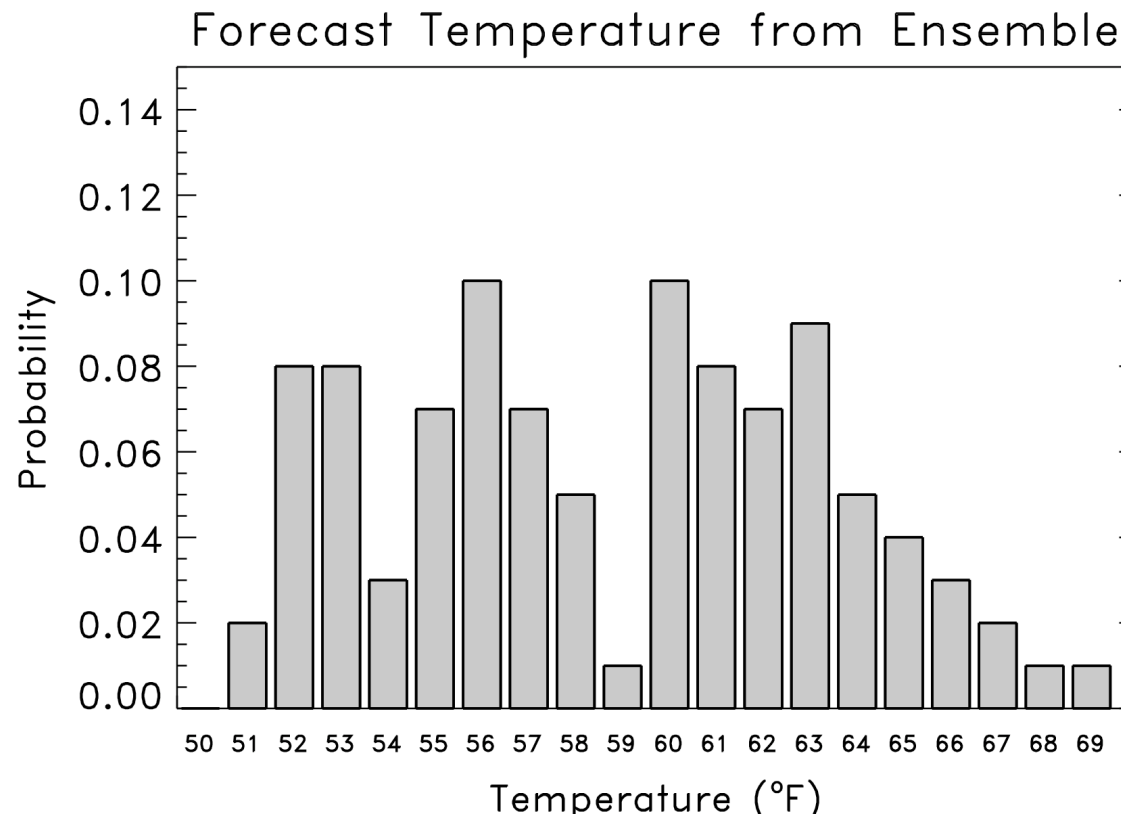$\beta$ = scale parameter

$$\Gamma(\alpha) = \int_0^\infty t^{\alpha-1}e^{-t}\,dt$$

### Forecast Probability Density

$\alpha = 1.0,\ \beta = 4.0$

$\alpha = 3.0\ \beta = 4.0$

Probability Density

Precipitation (mm)

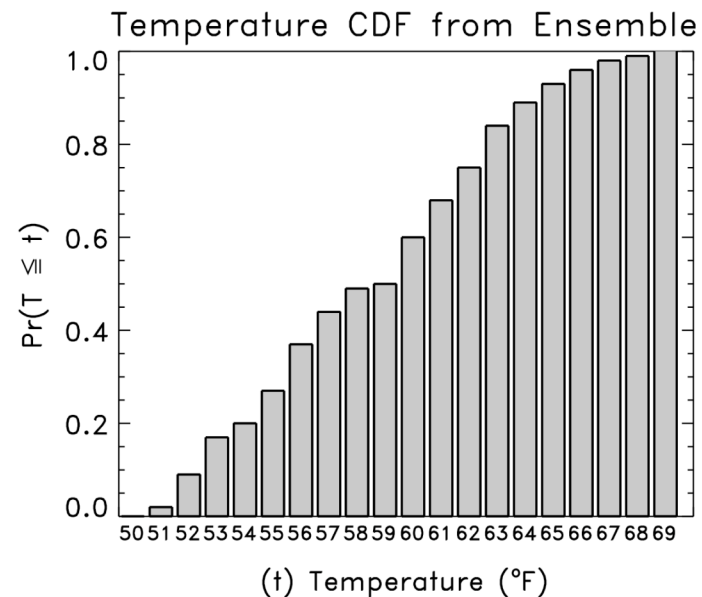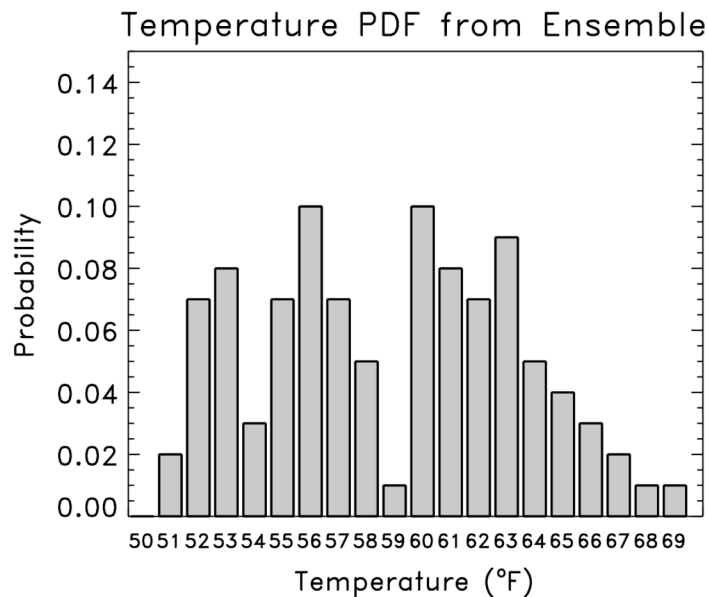# "Empirical" probability distributions

distribution derived from the data itself

# Cumulative Distribution Function (CDF)

- $F(t) = \Pr\{T \leq t\}$

  where T is the random variable, t is some specified threshold.



Temperature PDF from Ensemble

Temperature CDF from Ensemble

# Statistics

- **Definition**: "the science pertaining to the collection, analysis, interpretation or explanation, and presentation of data."

- **Goal**: make sure you understand terminology that we'll be using (mean, standard deviation, correlation, covariance, etc.)

# Measures of "location"

- T = [50, 51, 53, 54, 54,  57, 59, 63, 65, 66, 84] ($n$=11)

- Measure the centrality of this data set in some fashion.

- Mean (also called average, or 1st moment); minimizes RMS error:

$$\bar{T} = \frac{\sum\limits_{i=1}^{n} T_i}{n} = 59.63$$

- Median: central value of the sample, here = 57. Less affected by the 84 "outlier."  Minimizes mean absolute error.

# Measures of spread

- T = [50, 51, 53, 54, 54,  57, 59, 63, 65, 66, 84]
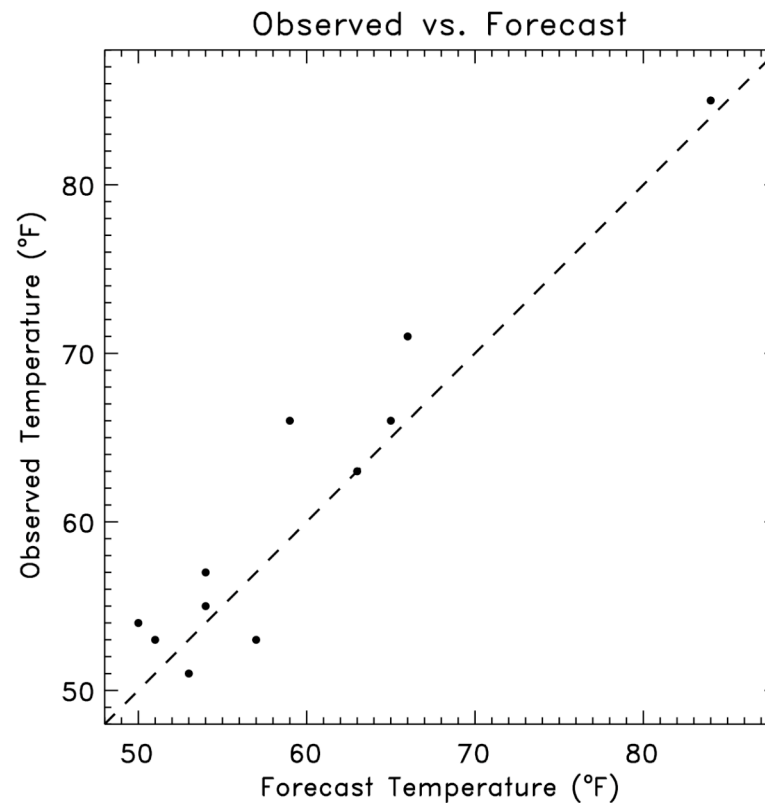
- Standard Deviation of sample:

  (variance is the square of this)

$$s = \left[ \frac{\sum_{i=1}^{n} \left(T_i - \bar{T}\right)^2}{n-1} \right]^{1/2} = 9.78$$

- IQR (Interquartile Range) = $q_{0.75}$ - $q_{0.25}$ = 65 - 53 = 12
  where $q_{0.75}$ is the 75th percentile (quantile) of the distribution and $q_{0.25}$ is the 25th percentile.

# Measures of association

- $T_f$ = [50, 51, 53, 54, 54, 57, 59, 63, 65, 66, 84]
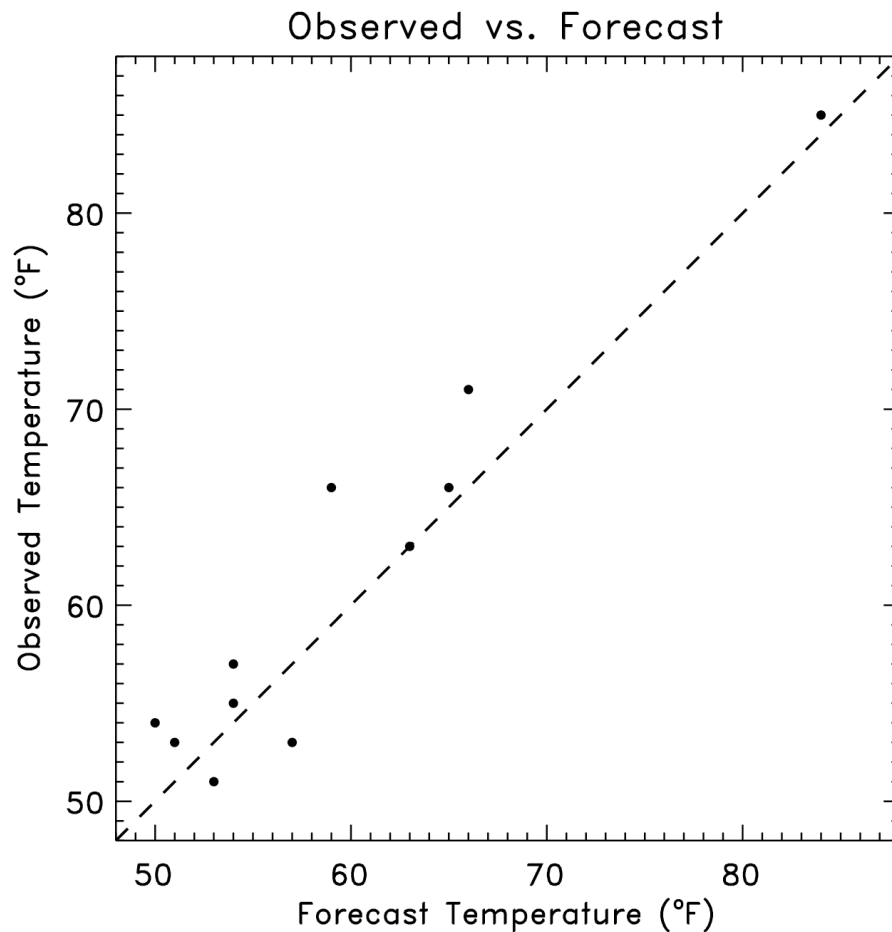- $T_o$ = [54, 53, 51, 57, 55, 53, 66, 63, 66, 71, 87]



Observed vs. Forecast

# Measures of association

- Pearson (ordinary) correlation:

$$r_{xy} = \frac{Cov(x,y)}{s_x s_y} = \frac{\dfrac{1}{n-1}\sum\limits_{i=1}^{n}\left[(x_i - \bar{x})(y_i - \bar{y})\right]}{\left\{\dfrac{1}{n-1}\sum\limits_{i=1}^{n}\left[(x_i - \bar{x})^2\right]^{1/2}\right\}\left\{\dfrac{1}{n-1}\sum\limits_{i=1}^{n}\left[(y_i - \bar{y})^2\right]^{1/2}\right\}}$$

$$= \frac{\sum\limits_{i=1}^{n}\left[x_i' \, y_i'\right]}{\left(\sum\limits_{i=1}^{n}\left[x_i'\right]^2\right)^{1/2}\left(\sum\limits_{i=1}^{n}\left[y_i'\right]^2\right)^{1/2}}$$

# Correlation, mean, standard deviation



Observed vs. Forecast

$r = 0.953$

$\overline{T}_f = 59.63$

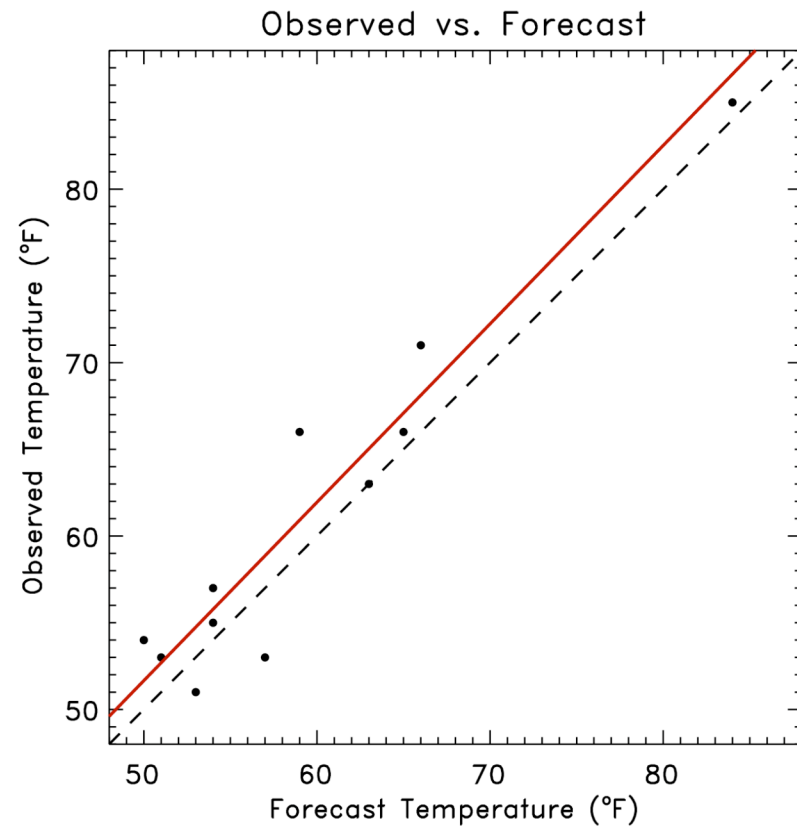$s\left(T_f\right) = 9.78$

$\overline{T}_o = 61.27$

$s\left(T_o\right) = 10.28$

# Regression

Find the equation that minimizes the squared difference between forecasts and observations.

$$T_o^{pred} = 1.478 + 1.006 * T^f$$

Methods like this used to statistically adjust weather forecasts.



Observed vs. Forecast

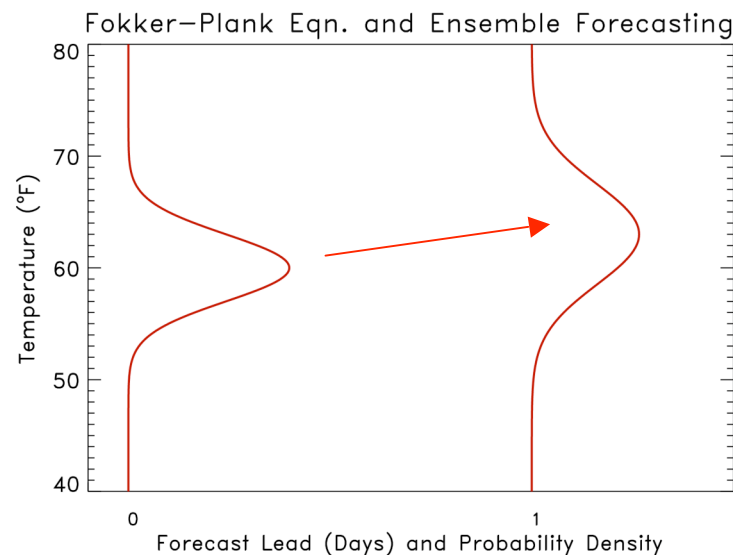# Connection between ensemble forecasts and PDFs

(there is a theory behind ensemble forecasting!)

Fokker-Planck equation to describe evolution of forecast PDF

$$\frac{\partial P\left(\mathbf{x}_t^t\right)}{\partial t} = -\nabla \bullet \left[ M\left(\mathbf{x}_t^t\right) P\left(\mathbf{x}_t^t\right) \right] + \sum_{i,j} \frac{\partial^2}{\partial \mathbf{x}_{t(i)}^t \partial \mathbf{x}_{t(j)}^t} \left( \frac{G\mathbf{Q}_t G^T}{2} \right)_{i,j} P\left(\mathbf{x}_t^t\right)$$

↑                                    ↑

errors due to chaos          errors due to the model



Fokker–Plank Eqn. and Ensemble Forecasting

Temperature (°F) vs Forecast Lead (Days) and Probability Density

(in reality, we never can get the pdf shown on day 1)

# Connection, cont'd



Fokker–Plank Eqn. and Ensemble Forecasting

In ensemble forecasting (ideally), we sample the initial pdf, and…

# Connection, cont'd



In ensemble forecasting (ideally), we sample the initial PDF, and evolve each initial condition forward with the forecast model(s) to randomly sample the day-1 PDF

# Questions?

# Baye's Rule

$$\Pr(A \ and \ E_1) = \Pr(A|E_1)\Pr(E_1)$$

$$= \Pr(E_1|A)\Pr(A)$$

combine 2 right-hand sides and rearrange

$$\Pr(E_1|A) = \frac{\Pr(A|E_1)\Pr(E_1)}{\Pr(A)} = \frac{\Pr(A|E_1)\Pr(E_1)}{\sum_{j=1}^{J}\Pr(A|E_j)\Pr(E_j)}$$